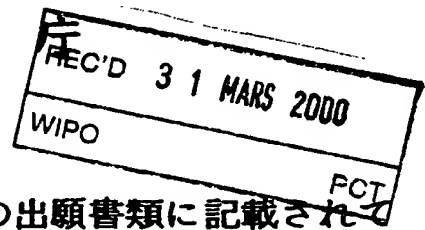


日本国特許
PATENT OFFICE
JAPANESE GOVERNMENT



別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office.

出願年月日

Date of Application:

1999年 3月 1日

出願番号

Application Number:

平成11年特許願第052156号

出願人

Applicant(s):

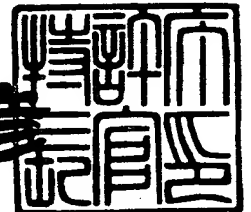
松下電器産業株式会社

PRIORITY
DOCUMENT
SUBMITTED OR TRANSMITTED IN
COMPLIANCE WITH RULE 17.1(a) OR (b)

2000年 1月14日

特許庁長官
Commissioner,
Patent Office

近藤 隆彦



出証番号 出証特平11-3093600

【書類名】 特許願

【整理番号】 2036610001

【提出日】 平成11年 3月 1日

【あて先】 特許庁長官

【国際特許分類】 G06K 09/20

【発明者】

 【住所又は居所】 大阪府門真市大字門真1006番地 松下電器産業株式会社内

 【氏名】 物部 祐亮

【発明者】

 【住所又は居所】 大阪府門真市大字門真1006番地 松下電器産業株式会社内

 【氏名】 広瀬 篤嗣

【発明者】

 【住所又は居所】 大阪府門真市大字門真1006番地 松下電器産業株式会社内

 【氏名】 梅林 明人

【特許出願人】

 【識別番号】 000005821

 【氏名又は名称】 松下電器産業株式会社

【代理人】

 【識別番号】 100083172

 【弁理士】

 【氏名又は名称】 福井 豊明

【手数料の表示】

 【予納台帳番号】 009483

 【納付金額】 21,000円

【提出物件の目録】

 【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

【包括委任状番号】 9713946

【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 文書画像処理装置及びそのタイトル領域抽出方法

【特許請求の範囲】

【請求項 1】 文書画像を複数の領域に分割する領域分割手段と、該領域分割手段によって分割された各領域について領域平均文字サイズを算出した後、該領域平均文字サイズに基づいて全領域の中からタイトル領域を抽出するタイトル領域抽出手段とを備えた文書画像処理装置において、

全領域の文字の平均高さに相当する全平均文字サイズを算出した後、該全平均文字サイズに抽出パラメータを乗算した抽出判定値と上記領域平均文字サイズとを比較し、上記抽出判定値より大きい領域平均文字サイズの領域をタイトル領域として抽出する上記タイトル領域抽出手段を備えたことを特徴とする文書画像処理装置。

【請求項 2】 上記タイトル領域抽出手段が、複数段階の抽出パラメータを用いて複数段階の上記抽出判定値を算出する請求項 1 に記載の文書画像処理装置。

【請求項 3】 上記タイトル領域抽出手段が、複数段階の抽出パラメータを用いて複数段階の上記抽出判定値を算出するとともに、抽出した段階を示すレベル属性を対応付けてタイトル領域を抽出する請求項 1 に記載の文書画像処理装置。

【請求項 4】 上記タイトル領域抽出手段が、領域平均文字サイズの最大値を全平均文字サイズで除算した値に基づいて上記複数段階の抽出パラメータを決定する請求項 2 または 3 に記載の文書画像処理装置。

【請求項 5】 上記タイトル領域抽出手段が、全平均文字サイズおよび領域平均文字サイズを、所定割合より大きい文字および所定割合より小さい文字を除外した文字より算出するトリム平均を用いる請求項 1 に記載の文書画像処理装置。

【請求項 6】 文書画像を複数の領域に分割し、各領域について文字の平均高さに相当する領域平均文字サイズを算出した後、該領域平均文字サイズに基づいて全領域の中からタイトル領域を抽出する文書画像処理装置のタイトル領域抽

出方法において、

全領域の全平均文字サイズを算出した後、該全平均文字サイズに抽出パラメータを乗算した抽出判定値と上記領域平均文字サイズとを比較し、上記抽出判定値より大きい領域平均文字サイズの領域をタイトル領域として抽出することを特徴とする文書画像処理装置のタイトル領域抽出方法。

【請求項 7】 複数段階の抽出パラメータを用いて複数段階の上記抽出判定値を算出する請求項 6 に記載の文書画像処理装置のタイトル領域抽出方法。

【請求項 8】 複数段階の抽出パラメータを用いて複数段階の上記抽出判定値を算出するとともに、抽出した段階を示すレベル属性を対応付けてタイトル領域を抽出する請求項 6 に記載の文書画像処理装置のタイトル領域抽出方法。

【請求項 9】 領域平均文字サイズの最大値を全平均文字サイズで除算した値に基づいて上記複数段階の抽出パラメータを決定する請求項 7 または 8 に記載の文書画像処理装置のタイトル領域抽出方法。

【請求項 10】 所定割合より大きい文字および所定割合より小さい文字を除外した文字の平均値を算出するトリム平均を用いて全平均文字サイズおよび領域平均文字サイズを算出する請求項 6 に記載の文書画像処理装置のタイトル領域抽出方法。

【請求項 11】 上記文書画像が複数頁の文書画像である請求項 6 ～ 10 に記載の文書画像処理装置のタイトル領域抽出方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、文書画像を画像データとして記憶・管理する文書画像処理装置と文書画像処理方法に関し、特に、上記文書画像からタイトル領域を抽出する上記装置と方法に関するものである。

【0002】

【従来の技術】

データ記憶装置の容量が著しく増加してきたことに伴って、スキャナ等から読み込んだ紙文書を画像データである文書画像として記憶・管理する文書画像処理

装置が急速に普及してきている。

【0003】

このような文書画像処理装置では、データ記憶装置に記憶された複数の文書画像の中から所望の文書画像を検索できるようにするため各文書画像にキーワードとなる文字列を対応付けて登録するようにしており、このキーワードは、文書画像を記憶する際にユーザが手入力していた。しかしながら、このようなキーワード入力作業をユーザが行うことは、文書数が多くなると作業量が膨大になるため現実的でない。そこで、近年では、文書画像に対して文字認識を行い、この認識結果である文字列をキーワードとすることによって、人手を介さずにキーワードを付与できるようにした装置も出現している。

【0004】

ところで、文書画像を効率良く検索するためには、適切なキーワードを付与しておくことが重要である。すなわち、ディスプレイに一覧表示された複数キーワードの中から所望の文書画像に対応するキーワードを特定するのが一般的な検索形態であるが、このようなキーワードを迅速に特定するためには、各キーワードが文書の内容を端的に表した内容でなければならない。

【0005】

特開平8-202859号公報では、タイトル文字列が属する領域（以下「タイトル領域」という。）を文書画像から抽出した後、このタイトル領域画像に対して文字認識を行い、この認識結果であるタイトル文字列をキーワードとする方法を提案している。タイトル文字列は文書の内容を端的に表した内容であるため、このようなタイトル領域抽出方法を採用した文書画像処理装置によれば所望の文書画像に対応するキーワードを迅速に特定できる。

【0006】

【発明が解決しようとする課題】

上記従来のタイトル領域抽出方法では、タイトル文字が当該文書画像に属する全ての文字の中で最も大きいサイズであるという観点から、文書画像を複数の領域（隣接する文字矩形相互を統合した領域）に分割して各領域内の平均文字サイズを算出し、この平均文字サイズが最も大きい領域をタイトル領域として抽出す

るようにしている。従って、このようなタイトル領域抽出方法によって抽出されるタイトル領域の数は、1つの文書画像につき当然1つとなる。

【0007】

しかしながら、近似した内容の複数の文書が存在する場合、タイトルも近似した内容となるのが通常であるため、上記従来のタイトル領域抽出方法には、近似した内容の文書が多数存在する場合所望の文書画像に対応するキーワードを迅速に特定できないという問題があった。

【0008】

上記問題を回避するため紙文書を作成する段階で似た内容のタイトルを付けないようにしてもよいが、このような準備作業をユーザに要求することは好ましくない。

【0009】

本発明は上記従来の事情に基づいて提案されたものであって、1つの文書画像から複数のタイトル領域を抽出できるようにした文書画像処理装置及びそのタイトル領域抽出方法を提供することを目的とするものである。

【0010】

【課題を解決するための手段】

本発明は上記目的を達成するために以下の手段を採用している。すなわち、図1に示すように、文書画像を複数の領域に分割する領域分割手段3と、該領域分割手段3によって分割された各領域について領域平均文字サイズを算出した後、該領域平均文字サイズに基づいて全領域の中からタイトル領域を抽出するタイトル領域抽出手段4とを備えた文書画像処理装置を前提としている。

【0011】

ここで、上記タイトル領域抽出手段4は、全領域の文字の平均高さに相当する全平均文字サイズを算出した後、該全平均文字サイズに抽出パラメータを乗算した抽出判定値と上記領域平均文字サイズとを比較し、上記抽出判定値より大きい領域平均文字サイズの領域をタイトル領域として抽出するようにしている。このようにすれば、上記抽出判定値より大きい領域平均文字サイズの領域であればタイトル領域として抽出されるため、1つの文書画像から複数のタイトル領域を抽

出できることになる。

【0012】

また、上記タイトル領域抽出手段4が、複数段階の抽出パラメータを用いて複数段階の上記抽出判定値を算出するようにしてもよい。このようにすれば、複数段階の抽出判定値に基づいて抽出判定がなされることになるため、タイトル領域だけでなくサブタイトル領域（タイトル文字より若干小さなサイズの文字からなるサブタイトル文字列が属する領域）をも抽出できる。

【0013】

更に、上記タイトル領域抽出手段4が、領域平均文字サイズの最大値を全平均文字サイズで除算した値に基づいて上記複数段階の抽出パラメータを決定するようにしてもよい。抽出パラメータを固定値とするのではなく領域平均文字サイズの最大値等に基づいて算出した方が良好な抽出判定値が得られる。

【0014】

上記全平均文字サイズ、領域平均文字サイズを求めるについて、所定割合より大きい文字および所定割合より小さい文字を除外したトリム平均を用いるとにより精度を上げることができる。

【0015】

【発明の実施の形態】

以下に本発明の実施の形態を図面に従って詳細に説明する。

（第1の実施の形態）

図1は、本発明を適用した文書画像処理装置の概略機能ブロック図であり、以下、その構成を文書画像登録手順とともに説明する。

【0016】

まず、例えばスキャナ等の文書画像入力手段1が紙文書を光電変換して多値画像データである文書画像8aを得、該文書画像は画像処理手段11aで記憶に適した処理（例えば圧縮処理）がなされて記憶手段8の文書画像エリアAaに登録される。もちろん画像処理手段11aを設けなくて多値画像データのまま文書画像エリアAaに登録しておいてもよい。

【0017】

上記文書画像入力手段 1 よりの文書画像は上記画像処理手段 11a に入力されるとともに、画像処理手段 11b にも入力されここで 2 値画像データに変換されて画像メモリ 7 に格納される。このように画像メモリ 7 に文書画像が格納された状態で、文字矩形生成手段 2 は上記画像メモリ 7 に記憶された文書画像を参照して、以下のラベリング処理を行う。このラベリング処理とは、注目する黒画素（以下「注目画素」という。）の上、右上、右、右下、下、左下、左、左上の 8 方向に隣接する画素のうち黒画素について当該注目画素と同一のラベル値（識別情報）を与える処理である。すなわち、図 7 に示すように $W1 \cdot W2 \cdot W3 \cdot W4 \cdot W6 \cdot W7 \cdot W8 \cdot W9$ の 8 画素が注目画素 $W5$ に連結する場合、文字矩形生成手段 2 は、黒画素である $W2 \cdot W3 \cdot W8$ に注目画素 $W5$ と同一のラベル値を与える。このようなラベリング処理を行うことによって、文書画像内の黒画素連結成分（連続する黒画素）毎に同一ラベル値を与えることができる。

【0018】

次いで、文字矩形生成手段 2 は、上記のように同一ラベル値を与えた黒画素連結成分を切り出すことによって文字矩形を生成し、この文字矩形を領域分割手段 3 に渡す。ここで、文字矩形とは黒画素連結成分の外接矩形を意味する。尚、文字によっては 1 つの黒画素連結成分で構成されていない場合もあり、このことを考慮して、上記ラベリング処理を行う前に文書画像中の黒画素領域を膨張させる処理をしておくこともできる。すなわち、注目する黒画素に隣接する 8 個の画素を黒画素に変換するという処理であり、この処理を適切な回数（通常 2、3 回）だけ施すことにより、黒画素の領域が拡大され 1 つの文字内で分離していた黒画素連結成分を 1 つに結合することができる。このような処理を行った上で、上記ラベリング処理を行うことにより、上記文字矩形を正しく文字毎に生成することが可能となる。

【0019】

上記文字矩形生成手段 2 の処理が終わると領域分割手段 3 は、各文字矩形について近傍を調べ、相互に隣接する文字矩形を統合することによって文書画像の領域を分割する。例えば図 8 に示す文字矩形 $C1 \sim C12$ を受けた領域分割手段 3 は、文字矩形 $C1 \sim C4 \cdot C5 \sim C9 \cdot C10 \sim C12$ をそれぞれ統合すること

によって文書画像を領域 $E1 \cdot E2 \cdot E3$ に分割する。このような領域分割処理を行うことによって、文書画像の領域を文字列毎に分割することができる。なお、文字矩形が相互に隣接している状態であるのか、あるいは、行間であるのか等の区別は文字間、行間に関する適当な閾値を用いて判定するようにしている。

【0020】

以上の結果、文書画像内における全ての文字サイズ（後述する）・分割された領域数・各領域内の文字矩形の数などの情報が得られる。本発明では、分割された各領域に対して1から始まる通し番号を付すとともに各領域に属する文字矩形に対しても1から始まる通し番号を付すようにしており、以下、 n 番目の領域内の文字矩形数を Num Char_n 、 n 番目の領域内における m 番目の文字サイズを $\text{Size Char}_n, m$ と表す。

【0021】

ところで図9に示すように、文字矩形の幅 $W1 \sim W4$ および面積 $A1 \sim A4$ は同一ポイント数の文字フォントを使用している場合であっても文字の種類に依存して大きく変動するのに対して、文字矩形の高さ $H1 \sim H4$ はこのような変動が小さい。従って本発明では、文字フォントのポイント数が比較的正確に反映される“文字矩形の高さ”を上記文字サイズとして採用するようにしている。

【0022】

ここで、タイトル領域抽出手段4は、上記のように分割された全領域のうち所定の領域のみをタイトル領域として抽出する。以下、このタイトル領域抽出処理を図2に示すフローチャートに従って説明する。

【0023】

まず、タイトル領域抽出手段4は各領域について領域平均文字サイズを算出する（図2、ステップ1）。この領域平均文字サイズとは1領域に属する全ての文字サイズの平均値であり、 n 番目の領域における領域平均文字サイズ SizeReg_n は、当該領域に属する全ての文字サイズ $\text{Size Char}_n, m$ の加算値を当該領域内の文字数 Num Char_n で除算した値となる。この関係を次式に示す。

【0024】

【数1】

$$\text{SizeReg}_n = \sum \text{SizeChar}_n / \text{NumChar}_n$$

次いで、上記のように算出した各領域の領域平均文字サイズ SizeReg_n と領域内の文字数 NumChar_n とから、文書画像内の全平均文字サイズ SizeAll を次式によって算出する（図2、ステップ2）。

【0025】

【数2】

$$\text{SizeAll} = \sum (\text{SizeReg}_n \times \text{NumChar}_n) / \sum \text{NumChar}_n$$

なお、領域平均文字サイズ SizeReg_n および全平均文字サイズ SizeAll の算出方法は上記した方法に限定されるものではなく、例えば、後に説明するトリム平均（最小値側および最大値側から所定割合例えば10%のデータを除外したうえで平均値を算出する方法）を採用することもできる。

【0026】

ここで、タイトル領域抽出手段4は、以下に示す抽出判定式が成立するか否かに基づいてタイトル領域の抽出判定を行う。

【0027】

【数3】

$$\text{SizeReg}_n \geq \text{SizeAll} \times \alpha$$

すなわち、上記のように算出した全平均文字サイズ SizeAll に抽出パラメータ α を乗算した値（抽出判定値）と各領域の領域平均文字サイズ SizeReg_n とを比較し、この抽出判定式が成立する領域のみをタイトル領域として抽出する（図2、ステップ3→4→5）。なお、抽出パラメータ α は1.0より大きな定数とし、1.2程度の値とするのが好ましい。

【0028】

以上の手順を繰り返し全ての領域について抽出判定が行われると（図2、ステップ3で“NO”）、タイトル領域抽出処理を終了し、ここで抽出された各タイトル領域画像8bは記憶手段8のタイトルエリアAbに登録される。

【0029】

次いで、文字認識手段5は、上記のように抽出されたタイトル領域画像を文書画像から切り出し、このタイトル領域画像に対して文字認識を行うことによって

文字コード列であるタイトル文字列を得た後、このタイトル文字列を文書登録手段 6 に渡す。

【0030】

上記タイトル文字列を受けた文書登録手段 6 は、記憶手段 8 での文書画像 8 a の格納ポインタ、上記タイトル領域画像 8 b の格納ポインタ・上記タイトル文字列・文書画像内におけるタイトル領域の位置およびサイズからなる登録情報を記憶手段 8 上のテーブルエリア A c に形成された登録情報管理テーブル 8 c (図 5 参照) に登録する。ここで、上記文書画像 8 a の格納ポインタは上記記憶手段 8 の文書画像エリア A a より得られ、上記タイトル画像 8 b の格納ポインタは上記記憶手段 8 のタイトルエリア A b より得られ、更に、タイトル領域の位置とサイズは文字認識手段 5 より得られることになる。

【0031】

このように登録情報管理テーブル 8 c が生成されると、以降に、キーボードやポインティングデバイス等からなる指示入力手段 9 より文書画像の検索が指示入力されると、表示制御手段 10 は、上記のように記憶されたタイトル領域画像およびタイトル文字列を図示しないディスプレイにリスト表示する(図 10 (I) 参照)。

【0032】

そして、上記リスト表示から所望のタイトル(タイトル領域画像またはタイトル文字列)が選択されたときには、このタイトルに対応する文書画像をディスプレイに表示する。このとき、図 10 (II) に示すように、矩形枠 F で囲むなどして文書画像内におけるタイトル領域を明示するのが好ましい。このような矩形枠 F は、登録情報管理テーブル 8 c に登録されているタイトル領域の位置およびサイズに基づいて生成できる。

【0033】

また、上記のようにディスプレイに表示されたリストからいずれか 1 つを選択する方法に加えて、指示入力手段 9 より特定のキーワードを入力し、該キーワードに該当するタイトルが登録情報管理テーブル 8 c に登録されているとき、対応する文書画像を表示するようにしてもよいことはもちろんである。

【0034】

以上のように本実施の形態によれば、抽出判定値より大きい領域平均文字サイズの領域であればタイトル領域として抽出する構成としているため、1つの文書画像から複数のタイトル領域を抽出できる。従って、似た内容の文書が多数存在する場合であっても、所望の文書画像に対応するキーワード（タイトル）を迅速に特定できる。

【0035】

なお、上記の説明では、タイトル領域抽出処理において抽出判定式の成立する領域が存在しなかった場合の手順については言及していないが、このような場合には、タイトル領域が抽出されなかった旨をディスプレイ表示するとともにキーワードとなる文字列を入力するようユーザに対して要求し、この要求に対してユーザが文字列を入力すると、この文字列を当該文書画像のタイトル文字列として用いるようにしている。

(第2の実施の形態)

上記第1の実施の形態では、抽出判定値より大きい領域平均文字サイズの領域であれば、領域平均文字サイズの大小を区別することなく同様にタイトル領域として抽出する構成としている。従って、タイトル文字より若干小さなサイズの文字からなるサブタイトル文字列はリスト表示せずタイトル文字列のみをリスト表示する処理など、領域平均文字サイズの大小に基づいた適切な処理を行うことができない。本実施の形態では、複数段階の抽出パラメータを用いて複数段階の抽出判定値を算出するとともにレベル属性（抽出した段階を示す情報）と対応付けてタイトル領域を抽出する構成とすることによって上記した問題を解消しており、以下、その構成を第1の実施の形態と異なる点のみ説明する。

【0036】

上記第1の実施の形態と同様の手順で領域平均文字サイズ $SizeReg_n$ および全平均文字サイズ $SizeAll$ を算出したタイトル領域抽出手段4は、以下に示す複数段階の抽出判定式が成立するか否かに基づいて複数段階の抽出判定を行う。

【0037】

【数4】

$$\text{SizeReg}_n \geq \text{SizeAll} \times \alpha_p$$

上式における α_p は、 p 段階（レベル p ）の抽出パラメータであり、【数 5】の条件を満たすように値を設定しておく。例えば、5 段階の抽出判定を行う場合には、 $\alpha_1=1.5$ 、 $\alpha_2=1.3$ 、 $\alpha_3=1.2$ 、 $\alpha_4=1.15$ 、 $\alpha_5=1.1$ 程度とするのが好ましい。

【0038】

【数 5】

$$1.0 < \alpha_p < \dots < \alpha_3 < \alpha_2 < \alpha_1$$

図 3 に示すフローチャートを用いて説明すると、タイトル領域抽出手段 4 は、レベル 1 から順に全レベルの抽出判定を行い（図 3、ステップ 14 → 15 → 14）、全レベルにおいて抽出判定式が成立しなかった場合には、この領域をタイトル領域として抽出せず、次の領域について抽出判定を行う（図 3、ステップ 14 → 13 → 14 → 15）。一方、いずれかのレベルにおいて抽出判定式が成立した場合には、この領域を当該レベルのタイトル領域として（上記レベル属性を対応付けて）抽出した後、次の領域について抽出判定を行う（図 3、ステップ 15 → 16 → 13 → 14 → 15）。

【0039】

以上の手順を繰り返し全ての領域について抽出判定が行われると（図 3、ステップ 13 で“NO”）、タイトル領域抽出処理を終了する。

【0040】

なお、抽出判定式の成立する領域が存在しなかった場合ユーザが入力した文字列をタイトル文字列として用いる点は上記第 1 の実施の形態と同様であり、このタイトル文字列のレベル属性はレベル 1、全レベル数も 1 としている。

【0041】

図 6 は、本実施の形態における登録情報管理テーブル 8 c の説明図であり、上記第 1 の実施の形態において示した構成にレベル属性フィールドと全レベル数フィールドとを加えた構成としている。そして文書登録手段 6 は、例えば 5 段階の抽出判定においてレベル 1 で抽出された領域がある場合、この領域に対応する全レベル数フィールドには“5”を、レベル属性フィールドには“1”をそれぞれ

登録する。

【0042】

図11は、本実施の形態の検索時においてディスプレイに表示される内容を示す図であり、上段にリスト表示するタイトルのレベル属性を指示入力手段9より範囲指定できるようにしている。そして、表示制御手段10は、登録情報管理テーブル8cのレベル属性フィールドと全レベル数フィールドとを参照することによって上記のように指定された範囲内のタイトルのみをディスプレイにリスト表示する。

【0043】

以上のように本実施の形態によれば、複数段階の抽出パラメータを用いて複数段階の抽出判定値を算出するとともにレベル属性と対応付けてタイトル領域を抽出する構成としているため、サブタイトル文字列はリスト表示せずタイトル文字列のみをリスト表示する処理など領域平均文字サイズの大小に基づいて、異なる処理を行うことができる。

(第3の実施の形態)

上記第2の実施の形態では、複数段階の抽出パラメータを予め設定する（固定値とする）構成としているが、このような抽出パラメータは入力された文書画像の特性に応じて決定するのが好ましい。本実施の形態では、領域平均文字サイズの最大値を全平均文字サイズで除算した値に基づいて複数段階の抽出パラメータを決定する（図4、ステップ23参照）ようにしており、以下、その構成を第2の実施の形態と異なる点のみ説明する。

【0044】

上記第2の実施の形態と同様の手順で領域平均文字サイズ $SizeReg_n$ および全平均文字サイズ $SizeAll$ を算出したタイトル領域抽出手段4は、まず、領域平均文字サイズの最大値 $\max \{SizeReg_n\}$ を全平均文字サイズ $SizeAll$ で除算した値 α_1 を次式によって算出する。

【0045】

【数6】

$$\alpha_1 = \max \{SizeReg_n\} / SizeAll$$

次いで、タイトル領域抽出手段4は、上記のように算出した α_1 と当該抽出判定の全レベル数 P ($P \geq 1$)とから、各レベルの抽出パラメータ α_p を次式によって決定する。

【0046】

【数7】

$$\alpha_p = \alpha_1 - (p - 1) \times (\alpha_1 - 1) / P$$

例えば α_1 が1.5で5段階の抽出判定を行う場合、各レベルの抽出パラメータ $\alpha_1 \sim \alpha_5$ は以下ようになる。

【0047】

【数8】

$$\alpha_1 = 1.5 - (1 - 1) \times (1.5 - 1) / 5 = 1.5$$

$$\alpha_2 = 1.5 - (2 - 1) \times (1.5 - 1) / 5 = 1.4$$

$$\alpha_3 = 1.5 - (3 - 1) \times (1.5 - 1) / 5 = 1.3$$

$$\alpha_4 = 1.5 - (4 - 1) \times (1.5 - 1) / 5 = 1.2$$

$$\alpha_5 = 1.5 - (5 - 1) \times (1.5 - 1) / 5 = 1.1$$

このように【数7】によれば、上記のように算出した α_1 から1.0の間で等間隔になるように各レベルの抽出パラメータ α_p を決定することができる。

【0048】

以降の手順は、上記のように決定した抽出パラメータを用いて抽出判定を行う点を除いて第2の実施の形態と同様であるため説明を省略する。

【0049】

ただし上記した方法には、文書画像内にタイトル領域が存在しない場合、 α_1 が例えば1.03など1.0付近の値となるため本文の領域をタイトル領域として誤抽出してしまうという不具合がある。そこで本発明では、例えば1.05など所定値以下となる抽出パラメータは採用しないようにしている。

【0050】

また、各レベル間の抽出パラメータの差が例えば0.03など所定値以下となると、良好な抽出判定ができないため、上記抽出パラメータの差が上記所定値(0.03)となるように抽出パラメータの設定値を修正するようにしている。すなわち上

記の場合、 α_1 から順に0.03ずつ減算した値を各レベルの抽出パラメータとして設定する。

【0051】

以上の結果全レベル数Pが減少する場合もあるが、このような場合には、実際のレベル数（全レベル数Pから減少レベル数を減じた値）を全レベル数Pとして登録情報管理テーブル8cの全レベル数フィールドに設定する。

【0052】

以上のように本実施の形態によれば、抽出パラメータを固定値とするのではなく、入力された文書画像の特性に応じて決定する構成としているため良好な抽出判定を行うことができる。

（第4の実施の形態）

上記の各実施の形態においては、全平均文字サイズの算出に比較的大小の大きいタイトル領域の文字も算入され、また、サイズの小さいコンマ、ピリオド、句読点も算入されるので、精度が低くなる傾向がある。そこで、文書画像の全文字から、所定割合（例えば90%）より大きいサイズの文字と、所定割合（例えば10%）より小さいサイズの文字を除外した文字から全平均文字サイズを算出する、いわゆるトリム平均を利用する。更に、領域平均文字サイズを算出するときにも、同様の問題が発生するところから、領域平均文字サイズの算出についても上記トリム平均を用いることもできる。

【0053】

これによって、全平均文字サイズ、および領域平均文字サイズとも、ピリオド、コンマ、句読点を除外した文字サイズを求めることができ、より精度の高い値が得られることになる。

【0054】

ここで、上記各実施の形態では領域平均文字サイズより、全平均文字サイズを算出しているが、同じ方法をこのトリム平均を用いる場合に適用すると、領域毎にサイズの大きい文字と小さい文字を除外することになるため、全平均文字サイズの算出においてタイトル領域に含まれるすべての文字を除外することができない。従ってここでは、全平均文字サイズを算出するときに、あらためて文書画像

中の全文字を対象として処理を行っている。

【0055】

但し、このトリム平均を用いる方式を使用するにしても、上記抽出パラメータとして、実施の形態1の所定値、あるいは実施の形態2、3の段階値のいずれを用いてもよいことはもちろんである。

【0056】

なお、上記の各実施の形態の説明では、文書画像となる文書の枚数については言及していないが、紙文書の枚数は特に限定されるものではない。すなわち、1枚であっても複数枚であっても、各頁に同じ抽出パラメータを用いる限り同様の効果が得られる。特に、実施の形態2、3においては、複数頁に対して同じ抽出パラメータを用いることにより、論文データのように複数頁に渡る単一文書から、タイトル、サブタイトルを正しく抽出することができる。

【0057】

また、上記の説明では、文字矩形の高さを文字サイズとして採用することとしているが、文字矩形の幅・面積を文字サイズとして採用してもよい。

【0058】

尚、図1の説明において、記憶手段8の前段と画像メモリ7の前段に画像処理手段11a、11bを設けて、タイトル抽出用の文書画像は2値画像データを用い、記憶手段8の文書画像エリアAaに登録される文書画像データとして、圧縮画像あるいは多値画像データを用いることができるようになっている。これによって、上記のように抽出されたタイトルに基づく検索処理の結果得られた文書画像をカラーで表示する等の多様な表示方法が可能となる。

【0059】

【発明の効果】

以上のように本発明によれば、抽出判定値より大きい領域平均文字サイズの領域をタイトル領域として抽出するようにしているため、1つの文書画像から複数のタイトル領域を抽出できる。また、複数段階の抽出パラメータに基づいて複数段階の抽出判定をすることもでき、更に、この複数段階の抽出パラメータを入力された文書画像の特性に応じて決定できる。また、全平均文字サイズの算出、あ

るいは領域平均文字サイズの算出に、大きい方の所定割合と小さい方の所定割合に属する文字を除外して算出するトリム平均を用いると、より精度を上げることができる。

【図面の簡単な説明】

【図 1】

本発明の文書画像処理装置の概略機能ブロック図である。

【図 2】

第 1 の実施の形態におけるタイトル領域抽出処理のフローチャートである。

【図 3】

第 2 の実施の形態におけるタイトル領域抽出処理のフローチャートである。

【図 4】

第 3 の実施の形態におけるタイトル領域抽出処理のフローチャートである。

【図 5】

第 1 の実施の形態における登録情報管理テーブルの説明図である。

【図 6】

第 2 の実施の形態における登録情報管理テーブルの説明図である。

【図 7】

ラベリング処理の説明図である。

【図 8】

領域分割処理の説明図である。

【図 9】

文字矩形の高さ・幅・面積の関係を示す図である。

【図 10】

第 1 の実施の形態の検索時においてディスプレイに表示される内容を示す図である。

【図 11】

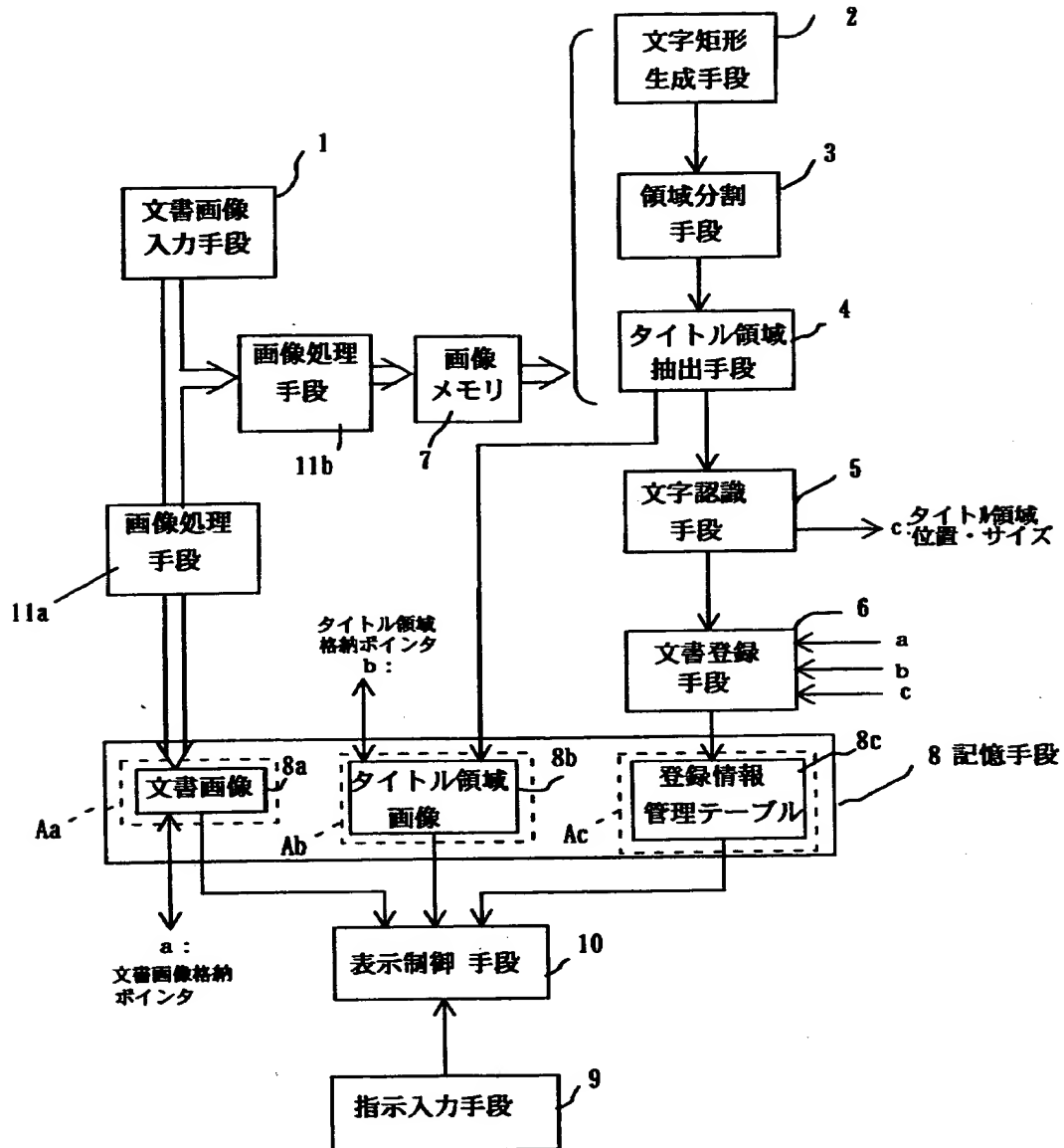
第 2 の実施の形態の検索時においてディスプレイに表示される内容を示す図である。

【符号の説明】

- 1 文書画像入力手段
- 2 文字矩形生成手段
- 3 領域分割手段
- 4 タイトル領域抽出手段
- 5 文字認識手段
- 6 文書登録手段
- 7 画像メモリ
- 8 記憶手段
- 9 指示入力手段
- 1 0 表示制御手段
- 1 1 (1 1 a、1 1 b) 画像処理手段

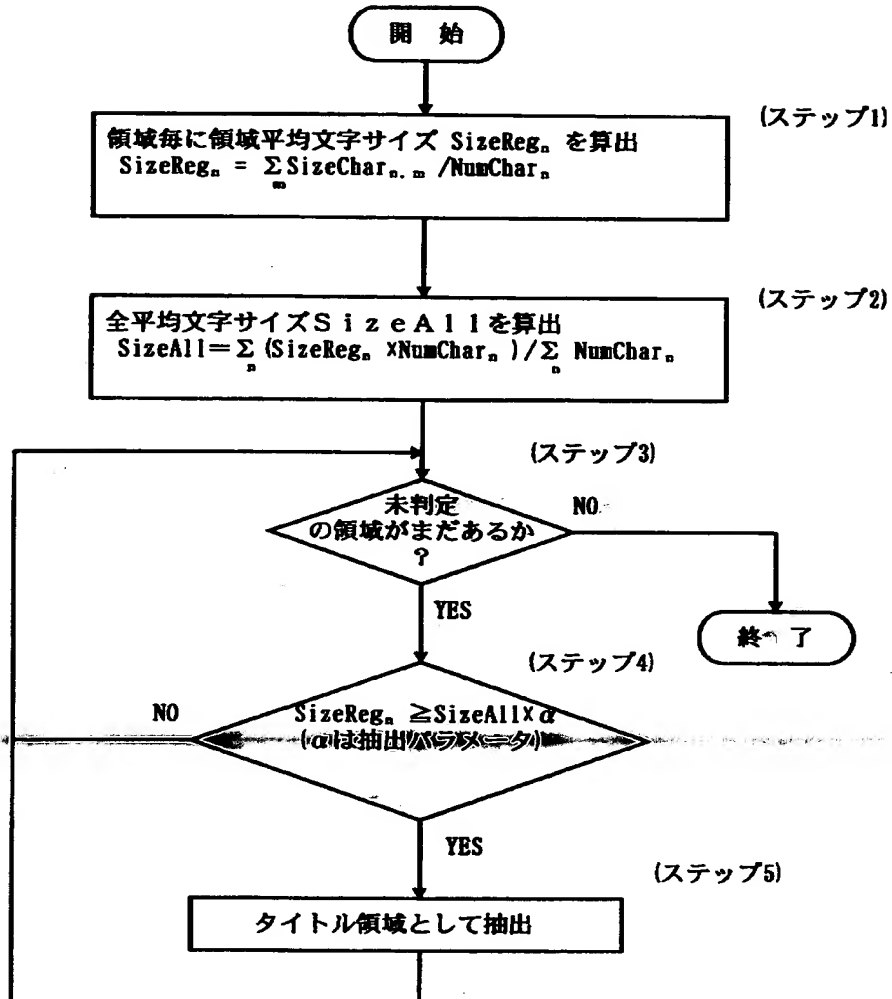
【書類名】 図面

【図 1】



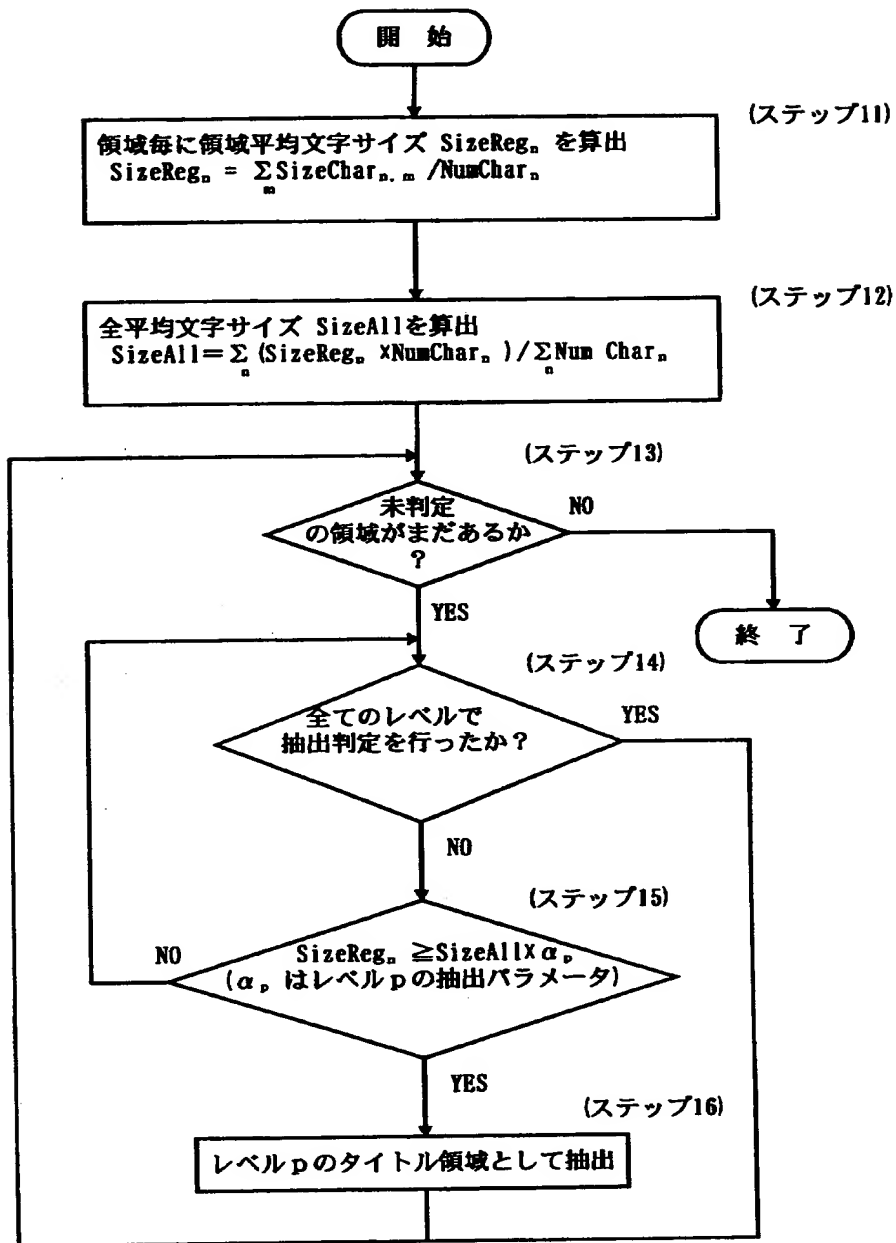
【図 2】

第 1 の実施の形態におけるタイトル領域抽出処理のフローチャート



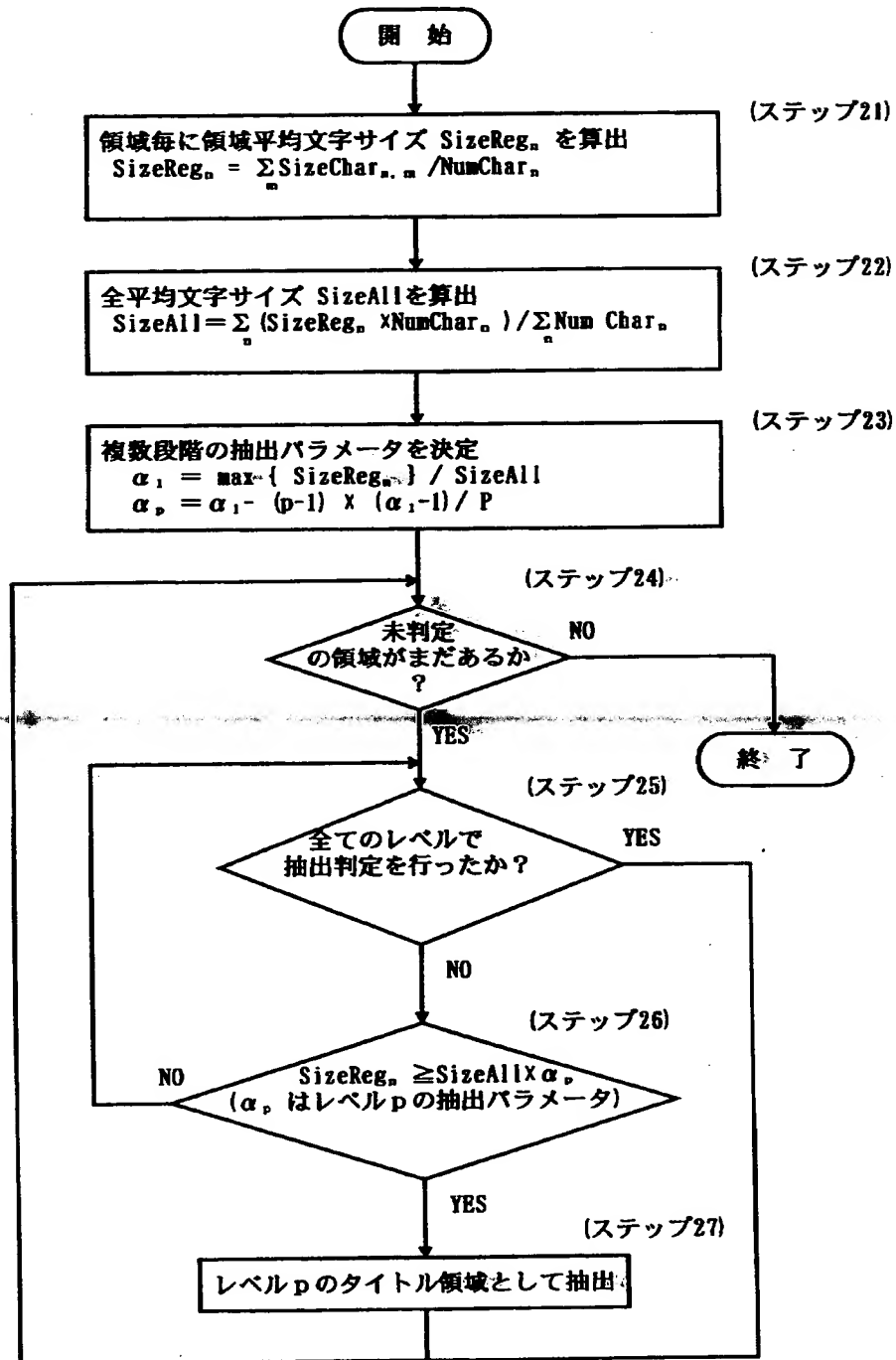
【図 3】

第 2 の実施の形態におけるタイトル領域抽出処理のフローチャート



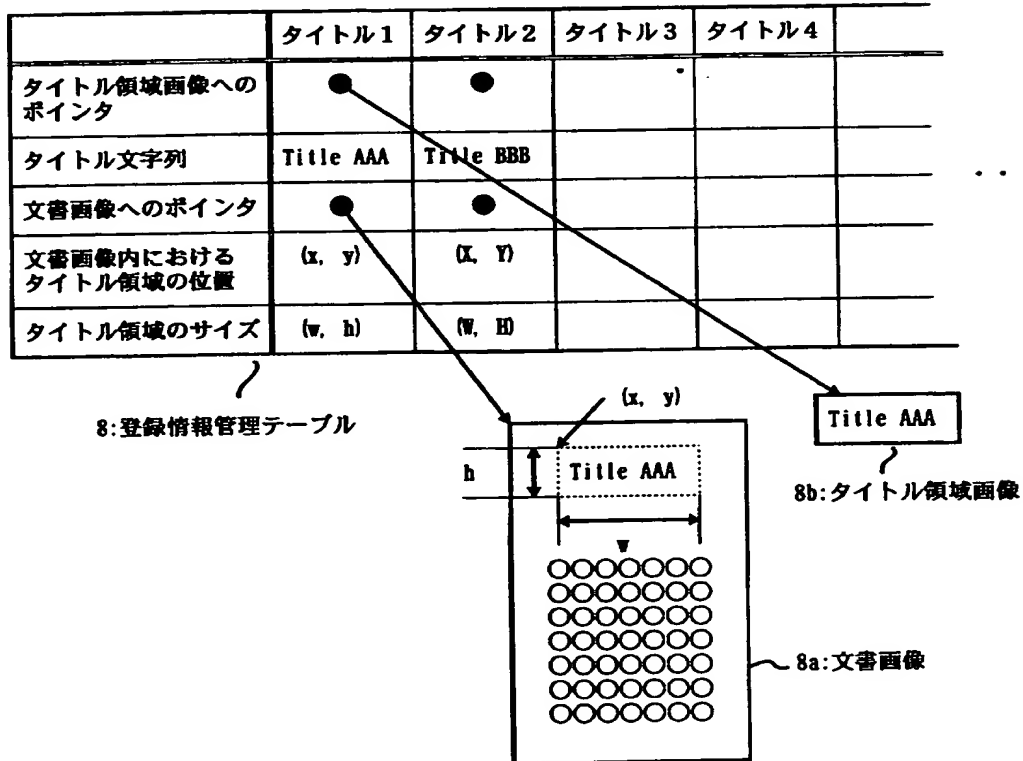
【図 4】

第 3 の実施の形態におけるタイトル領域抽出処理のフローチャート



【図 5】

第 1 の実施の形態における登録情報管理テーブルの説明図

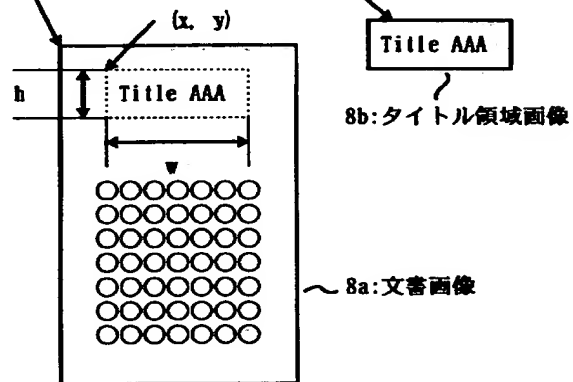


【図 6】

第 2 の実施の形態における登録情報管理テーブルの説明図

| | タイトル 1 | タイトル 2 | タイトル 3 | タイトル 4 | |
|------------------------|-----------|-----------|--------|--------|-----|
| タイトル領域画像への ポインタ | ● | ● | | | |
| タイトル文字列 | Title AAA | Title BBB | | | |
| 文書画像へのポインタ | ● | ● | | | |
| 文書画像内における タイトル領域の位置 | (x, y) | (X, Y) | | | ... |
| タイトル領域のサイズ | (w, h) | (W, H) | | | |
| レベル属性 | 1 | 3 | | | |
| 全レベル数 | 5 | 5 | | | |

8:登録情報管理テーブル



8b:タイトル領域画像

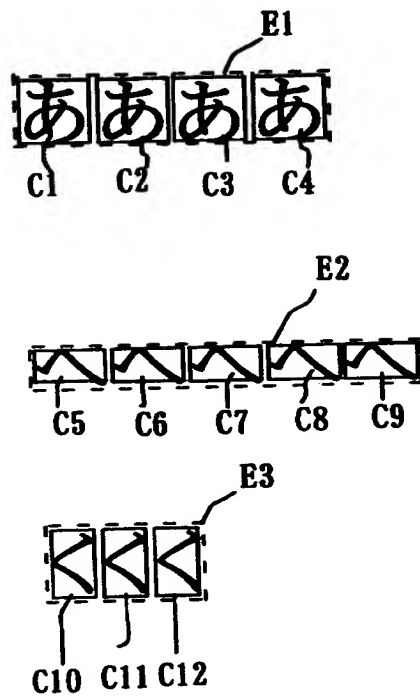
8a:文書画像

【図 7】

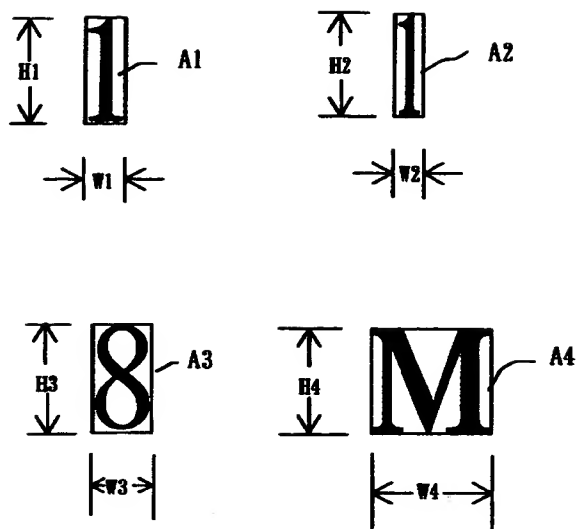
ラベリング処理の説明図

| | | |
|------------|------------|------------|
| W 1 (白) | W 2 (黒) | W 3 (黒) |
| W 4 (白) | W 5 (黒) | W 6 (白) |
| W 7 (白) | W 8 (黒) | W 9 (白) |

【図 8】



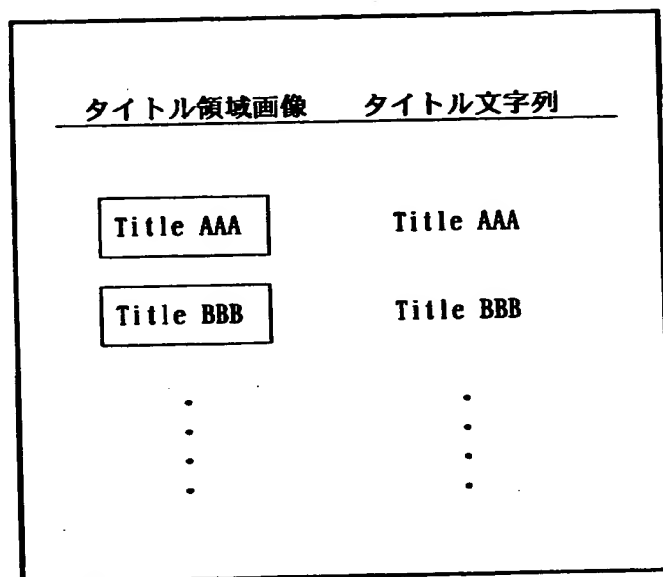
【図9】



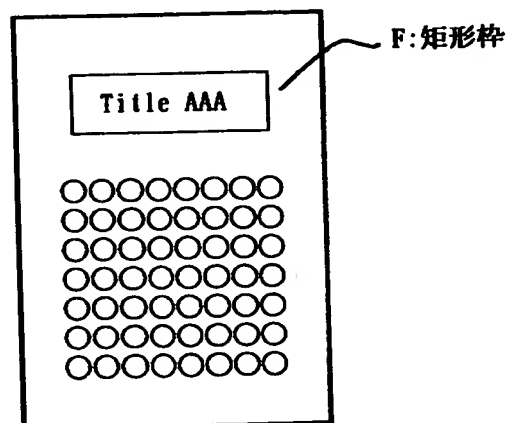
【図 10】

第 1 の実施の形態の検索時においてディスプレイに表示される内容

(I)



(II)



【図 11】

第 2 の実施の形態の検索時においてディスプレイに表示される内容

| リスト表示するレベル属性範囲の指定 | | |
|-------------------|-----------|---------------|
| レベル | 1 | 3 |
| タイトル領域画像 | タイトル文字列 | レベル属性 / 全レベル数 |
| Title AAA | Title AAA | 1 / 5 |
| Title BBB | Title BBB | 3 / 5 |
| ⋮ | ⋮ | ⋮ |
| ⋮ | ⋮ | ⋮ |
| ⋮ | ⋮ | ⋮ |
| ⋮ | ⋮ | ⋮ |

【書類名】 要約書

【要約】

【課題】 文書画像よりタイトルを抽出する処理の精度を上げる。

【解決手段】 上記タイトル領域抽出手段4は、全領域における文字の平均高さに相当する全平均文字サイズを算出した後、該全平均文字サイズに抽出パラメータを乗算した抽出判定値と領域平均文字サイズとを比較し、上記抽出判定値より大きい領域平均文字サイズの領域をタイトル領域として抽出するようにしている。このようにすれば、上記抽出判定値より大きい領域平均文字サイズの領域であればタイトル領域として抽出されるため、1つの文書画像から複数のタイトル領域を抽出できることになる。

【選択図】 図 1

出 願 人 履 歴 情 報

識別番号

[000005821]

1. 変更年月日

1990年 8月28日

[変更理由]

新規登録

住 所

大阪府門真市大字門真1006番地

氏 名

松下電器産業株式会社